Object Detection

Object Detection

- The goal of object detection is to localize objects in an image and tell their class
- Localization: place a tight bounding box around object
- Most approaches find only objects of one or a few specific classes, e.g. car or cow





Type of Approaches

Different approaches tackle detection differently. They can roughly be categorized into three main types:

• Find interest points, followed by Hough voting

- Compute interest points (e.g., Harris corner detector is a popular choice)
- Vote for where the object could be given the content around interest points



- Compute interest points (e.g., Harris corner detector is a popular choice)
- Vote for where the object could be given the content around interest points



- Compute interest points (e.g., Harris corner detector is a popular choice)
- Vote for where the object could be given the content around interest points



- Compute interest points (e.g., Harris corner detector is a popular choice)
- Vote for where the object could be given the content around interest points



- Compute interest points (e.g., Harris corner detector is a popular choice)
- Vote for where the object could be given the content around interest points



Type of Approaches

Different approaches tackle detection differently. They can roughly be categorized into three main types:

- Find interest points, followed by Hough voting
- **Sliding windows**: "slide" a box around image and classify each image crop inside a box (contains object or not?)

• Slide window and ask a classifier: "Is sheep in window or not?"



0.1 confidence

• Slide window and ask a classifier: "Is sheep in window or not?"



-0.2

• Slide window and ask a classifier: "Is sheep in window or not?"



-0.1

• Slide window and ask a classifier: "Is sheep in window or not?"



• Slide window and ask a classifier: "Is sheep in window or not?"



. . . 1.5

• Slide window and ask a classifier: "Is sheep in window or not?"



• Slide window and ask a classifier: "Is sheep in window or not?"



• Slide window and ask a classifier: "Is sheep in window or not?"



• Slide window and ask a classifier: "Is sheep in window or not?"



0.1 confidence-0.2 -0.1 0.1 ... 1.5 ... 0.5 0.4 0.3

[Slide: R. Urtasun]

Type of Approaches

Different approaches tackle detection differently. They can roughly be categorized into three main types:

- Find interest points, followed by Hough voting
- Sliding windows: "slide" a box around image and classify each image crop inside a box (contains object or not?)
- Generate region (object) proposals, and classify each region

• Group pixels into object-like regions



• Group pixels into object-like regions



• Group pixels into object-like regions



• Generate many different regions



• Generate many different regions



• Generate many different regions



• The hope is that at least a few will cover real objects



• The hope is that at least a few will cover real objects



• Select a region



• Crop out an image patch around it, throw to classifier (e.g., Neural Net)



classifier ``dog" or not?

confidence: -2.5

• Do this for every region



• Do this for every region



• Do this for every region





Type of Approaches

Different approaches tackle detection differently. They can roughly be categorized into three main types:

- **Sliding windows**: "slide" a box around image and classify each image crop inside a box (contains object or not?)
- Generate region (object) proposals, and classify each region

Object Detection via Hough Voting: Implicit Shape Model

B. Leibe, A. Leonardis, B. Schiele

Robust Object Detection with Interleaved Categorization and

Segmentation

IJCV, 2008

Paper: http://www.vision.rwth-aachen.de/publications/pdf/leibe-interleaved-ijcv07final.pdf

Start with Simple: Line Detection

• How can I find lines in this image?



[Source: K. Grauman]

Sanja Fidler

Hough Transform

- Idea: Voting (Hough Transform)
- Voting is a general technique where we let the features vote for all models that are compatible with it.
 - Cycle through features, cast votes for model parameters.
 - Look for model parameters that receive a lot of votes.

[Source: K. Grauman]
• Hough space: parameter space



- Connection between image (x, y) and Hough (m, b) spaces
 - A line in the image corresponds to a point in Hough space
 - What does a point (x_0, y_0) in the image space map to in Hough space?

[Source: S. Seitz]

• Hough space: parameter space



- Connection between image (x, y) and Hough (m, b) spaces
 - A line in the image corresponds to a point in Hough space
 - A point in image space votes for all the lines that go through this point. This votes are a line in the Hough space.
- [Source: S. Seitz]

• Hough space: parameter space



Two points: Each point corresponds to a line in the Hough space
A point where these two lines meet defines a line in the image!
[Source: S. Seitz]

• Hough space: parameter space



- Vote with each image point
- Find peaks in Hough space. Each peak is a line in the image.

[Source: S. Seitz]

- Issues with usual (m, b) parameter space: undefined for vertical lines
- A better representation is a polar representation of lines



d: perpendicular distance from line to origin

 $\ensuremath{\boldsymbol{\theta}}$: angle the perpendicular makes with the x-axis

 $x\cos\theta - y\sin\theta = d$

Point in image space \rightarrow sinusoid segment in Hough space

[Source: S. Seitz]

Example Hough Transform

With the parameterization $x \cos \theta + y \sin \theta = d$

- Points in picture represent sinusoids in parameter space
- Points in parameter space represent lines in picture
- Example 0.6x + 0.4y = 2.4, Sinusoids intersect at d = 2.4, $\theta = 0.9273$



[Source: M. Kazhdan, slide credit: R. Urtasun]

• Hough Voting algorithm

Using the polar parameterization:

 $x\cos\theta - y\sin\theta = d$

Basic Hough transform algorithm

1. Initialize H[d, θ]=0 2. for each edge point I[x,y] in the image for θ = [θ_{min} to θ_{max}] // some quantization $d = x \cos \theta - y \sin \theta$ H[d, θ] += 1



H: accumulator array

3. Find the value(s) of (d, $\theta)$ where H[d, $\theta]$ is maximum

4. The detected line in the image is given by $d = x \cos\theta - y \sin\theta$

[Source: S. Seitz]

• What about circles? How can I fit circles around these coins?



Assume we are looking for a circle of known radius r

• Circle:
$$(x - a)^2 + (y - b)^2 = r^2$$

- Hough space (a, b): A point (x_0, y_0) maps to $(a - x_0)^2 + (b - y_0)^2 = r^2 \rightarrow a$ circle around (x_0, y_0) with radius r
- Each image point votes for a circle in Hough space



Each point in geometric space (left) generates a circle in parameter space (right). The circles in parameter space intersect at the (a, b) that is the center in geometric space.

[Source: H. Rhody]

What if we don't know r?

• Hough space: ?



[Source: K. Grauman]

What if we don't know r?

• Hough space: conics



[Source: K. Grauman]



[Source: K. Grauman]

Iris detection



Gradient+threshold

Hough space (fixed radius)

Max detections

[Source: K. Grauman]

Generalized Hough Voting

• Hough Voting for general shapes



Offline procedure:

At each boundary point, compute displacement vector: $\mathbf{r} = \mathbf{a} - \mathbf{p}_{i}$.

Store these vectors in a table indexed by gradient orientation θ .

Dana H. Ballard, Generalizing the Hough Transform to Detect Arbitrary Shapes, 1980

Implicit Shape Model

- Implicit Shape Model adopts the idea of voting
- Basic idea:
 - Find interest points in an image
 - Match patch around each interest point to a training patch
 - Vote for object center given that training instance

Vote for object center



vote for center of object

• Vote for object center



vote for center of object

• Vote for object center



vote for center of object

Vote for object center



of course some wrong votes are bound to happen...

Vote for object center



But that's ok. We want only **peaks** in voting space.

• Find the patches that produced the peak



Find patches that voted for the peaks (back-projection).

• Place a box around these patches \rightarrow objects!



Find full objects based on the back-projected patches.

• Really easy. Only one problem... Would be slow... How do we make it fast?



we need to match a patch around each yellow + to all patches in all training images $\ \rightarrow\$ SLOW

Sanja Fidler

CSC420: Intro to Image Understanding

- Visual vocabulary (we saw this for retrieval)
- Compare each patch to a small set of visual words (clusters)



Visual words (visual codebook)!

• Training: Getting the vocabulary

training image



• Find interest points in each training image

training image



detect interest points (e.g. Harris)

• Collect patches around each interest point

training image



extract an image patch around each interest point

• Collect patches across all training examples

training images



collect all patches



• Cluster the patches to get a small set of "representative" patches

training images



- cluster the patches to get a few ``representative'' patches
- each cluster represented as the average of all patches that belong to the cluster

collect all patches



visual codebook



Implicit Shape Model: Training

- Represent each training patch with the closest visual word.
- Record the displacement vectors for each word across all training examples.



Training image



Visual codeword with displacement vectors

[Leibe et al. IJCV 2008]

Implicit Shape Model: Test

- At test times detect interest points
- Assign each patch around interest point to closes visual word
- Vote with all displacement vectors for that word



[Source: B. Leibe]

Recognition Pipeline



[Source: B. Leibe]

Recognition Summary

- Apply interest points and extract features around selected locations.
- Match those to the codebook.
- Collect consistent configurations using Generalized Hough Transform.
- Each entry votes for a set of possible positions and scales in continuous space.
- Extract maxima in the continuous space using Mean Shift.
- Refinement can be done by sampling more local features.

[Source: R. Urtasun]



Original image

[Source: B. Leibe, credit: R. Urtasun]



Interest points

[Source: B. Leibe, credit: R. Urtasun]

Example



Matched patches

[Source: B. Leibe, credit: R. Urtasun]


Voting space



1st hypothesis



2nd hypothesis



3rd hypothesis

Scale Invariant Voting

Scale-invariant feature selection

- Scale-invariant interest points
- Rescale extracted patches
- Match to constant-size codebook

Generate scale votes

• Scale as 3rd dimension in voting space

• Search for maxima in 3D voting space [Source: B. Leibe, credit: R. Urtasun]

Scale Invariant Voting



[Slide credit: R. Urtasun]

Scale Voting: Efficient Computation

Continuous Generalized Hough Transform

- Binned accumulator array similar to standard Gen. Hough Transf.
- Quickly identify candidate maxima locations
- Refine locations by Mean-Shift search only around those points
- Avoid quantization effects by keeping exact vote locations.



[Source: B. Leibe, credit: R. Urtasun]

Extension: Rotation-Invariant Detection

- Polar instead of Cartesian voting scheme
- Recognize objects under image-plane rotations
- Possibility to share parts between articulations
- But also increases false positive detections





Sometimes it's Necessary



[Source: B. Leibe, credit: R. Urtasun]

Recognition and Segmentation



augment each cluster with a figure-ground mask

• Augment each visual word with meta-deta: for example, segmentation mask

Recognition and Segmentation



Results



[Source: B. Leibe]



[Source: B. Leibe]

Results



[Source: B. Leibe]

Results



[Source: B. Leibe]

Inferring Other Information: Part Labels



[Source: B. Leibe]

Sanja Fidler

45 / 48

Inferring Other Information: Part Labels



[Source: B. Leibe]

Inferring Other Information: Depth



"Depth from a single image"





[Source: B. Leibe]

Conclusion

- Exploits a lot of parts (as many as interest points)
- Very simple Voting scheme: Generalized Hough Transform
- Works well, but not as well as Deformable Part-based Models with latent SVM training (next time)
- Extensions: train the weights discriminatively.
- Code, datasets & several pre-trained detectors available at http://www.vision.ee.ethz.ch/bleibe/code